

A Modern Analysis of Hutchinson’s Trace Estimator.

1st Maciej Skorski
University of Luxembourg

Abstract—The paper establishes the new state-of-art in the accuracy analysis of Hutchinson’s trace estimator. Leveraging tools that have not been previously used in this context, particularly hypercontractive inequalities and concentration properties of sub-gamma distributions, we offer an elegant and modular analysis, as well as numerically superior bounds. Besides these improvements, this work aims to better popularize the aforementioned techniques within the CS community.

Index Terms—Randomized algorithms, Trace estimation, Hutchinson’s estimator, Monte Carlo methods, Hypercontractive inequalities, Sub-gamma distributions

I. INTRODUCTION

A. Background

Estimating the matrix trace is a problem of fundamental interest [1]–[3] and arises in many problems such as approximating spectral properties of matrices [4]–[6], solving partial differential equations [7]–[11], error evaluation in machine-learning [12], and combinatorial counting [13] to mention a few. For readers particularly interested in data science or optimization, it is of critical interest as the hessian trace stores valuable information about curvature. To give a meaningful example, consider fitting a linear classifier on the MNIST dataset [14] of hand-written digits, widely used as a benchmark and a toy problem in machine learning. Since images are in resolution 28×28 and grouped in 10 classes, the number of parameters is $m = 28 \cdot 28 \cdot 10 = 7840$, and the size of the hessian matrix is $m^2 \approx 6 \cdot 10^6$. The diagonal average, proportional to the trace, equals the average eigenvalue (valuable information), whereas individual rows or columns are nearly zero (up to random noise). This is illustrated in Figure 1.

As seen from the above example, already toy examples lead to huge matrices; for larger problems storing and inspecting such matrices is impossible. Fortunately, it is possible to estimate the trace (equal to the diagonal sum, or equivalently to the sum of eigenvalues) *without knowing the full matrix*. The requirement is that one can compute efficiently *matrix-vector* products; in the case of the hessian such products can be computed by by auto-differentiation [15] in all popular machine-learning frameworks such as Tensorflow [16], PyTorch [17] and JAX [18]. Under the assumption that an $m \times m$ matrix A is *symmetric positive semi-definite* (which applies to all Hessians), the popular estimator due to Hutchinson [1] is

$$\text{tr}_H(A) \triangleq z^T A z, \quad z \sim^{iid} \{-1, 1\}^m, \quad (1)$$

²https://github.com/maciejskorski/ml_examples/blob/master/TF_Hess_Visualize.ipynb

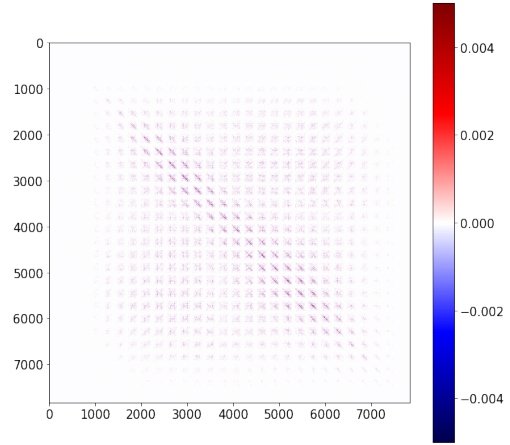


Fig. 1. Hessian matrix of the linear classifier trained on MNIST data. The diagonal stores valuable information as opposed to rows and columns: entries in rows and columns are nearly zero, whereas the diagonal average is ≈ 0.25 . The Tensorflow code is available as a Python Notebook at GitHub².

where the components of the vector z are *Rademacher random variables*, equal to $+1$ or -1 with equal probability. The above estimator uses only one sampled v and thus is noisy (of high-variance), so it is usually boosted by averaging n independent trials (for suitably large n). Formally:

$$\text{tr}_{H^{(n)}}(A) \triangleq \frac{1}{n} \sum_{i=1}^n z_i^T A z_i, \quad z_i \sim^{iid} \{-1, 1\}^m. \quad (2)$$

The estimators are unbiased, that is correct on average:

$$\mathbb{E} \text{tr}_H(A) = \mathbb{E} \text{tr}_{H^{(n)}}(A) = \text{tr}(A), \quad (3)$$

and this is fairly easy to prove. The focus of this work is on a more challenging problem of understanding their *concentration properties*. Here we ask the following question:

What sample size n *guarantees* the desired relative error ϵ at the confidence level of $1 - \delta$?

Formally, for the error $\text{err}_{H^{(n)}}(A) \triangleq \frac{\text{tr}_{H^{(n)}}(A)}{\text{tr}(A)} - 1$, and fixed ϵ, δ , we want possibly small $n = n(\epsilon, \delta)$ such that

$$|\text{err}_{H^{(n)}}(A)| \leq \epsilon \quad \text{with prob. at least } 1 - \delta. \quad (4)$$

B. Related Work

The estimator is already more than 30 years old [1], and although alternatives exist (such as methods based on Lanczos

quadrature [19]), it is provably best in terms of variance [20] and arguably wins in simplicity, being preferred by developers of statistical/learning software [21]. Quite surprisingly, it had been used for a while without a rigorous assessment of accuracy, until the work of Avron and Toledo [3], who established the finite sample size guarantees. Their result was later improved by Roosta-Khorasani and Ascher [20], who essentially got rid of the lossy proof step involving a crude union bound. Their approach is based on the Chernoff-like direct calculations, and offers the bound $n = O(\epsilon^{-2} \log(1/\delta))$.

II. CONTRIBUTION

A. Summary

In this work we offer a *modular, modern and more accurate* analysis of Hutchinson’s estimator. The improvements upon prior works can be summarized as follows:

- **Modularity and Novelty of Techniques.** In the first step, we explicitly link the estimator accuracy to the dispersion of the *Rademacher Quadratic Chaos* with a unitary matrix kernel. We then obtain a bound on such chaoses by means of *Hypercontractive Inequality*, an important result in Boolean Fourier Analysis [22], originally due to Aline Bonami [23]. As the final step we express this bound as the *sub-gamma property* (see the works of Boucheron [24]), which makes it convenient (and accurate) to conclude the final concentration result. Moreover, we state our results for the multiple-sample and one-shot estimators; this is of interest, since the base building block may be boosted differently than by averaging (see for example the median algorithm [25]). Thus, in contrast to ad hoc calculations in prior works, we are able to use established tools from the field of high-dimensional probability and Fourier analysis. Our transparent approach opens a way for further refinements.
- **Superior Accuracy.** With the dedicated tools we obtain bounds that are numerically better up to an order of magnitude. They are stated as elementary formulas, convenient to use in practical applications of trace estimation.

B. Preliminaries

Below we explain the notation and definitions used when formulating our results presented in the next section.

The d -th norm of a r.v. X is $\|X\|_d \triangleq (\mathbb{E}|X|^d)^{1/d}$; it is a valid norm (e.g. it obeys the triangle inequality) when $d \geq 1$ [25]. A useful tool for studying concentration is the *moment generating function*, defined as $\text{MGF}_X(t) = \mathbb{E} \exp(tX)$ (a function of real parameter t). In modern high-dimensional probability one classifies distributions according to the behaviour of MGF; in particular we call X *sub-gamma* with variance factor v and scale c when $\log \text{MGF}_{|X|}(t) \leq \frac{vt^2}{2(1-ct)}$, and denote by $X \in \Gamma(v, c)$ [24]; the same formula holds for the centered gamma distribution, hence the name.

C. Results

In what follows we assume that A is non-zero symmetric positive semi-definite matrix. Then, necessarily, $\text{tr}(A) > 0$.

1) *One-Shot Estimator:* Below we present the following bound on the accuracy of the estimator in (1):

Theorem 1 (Hutchison’s Estimator 1-Sample Bound). *For any integer $d \geq 2$, the relative error of Hutchison’s Estimator (1) obeys the following bound*

$$\|\text{err}_H(A)\|_d \leq d - 1. \quad (5)$$

In particular, this implies the sub-gamma behaviour

$$\text{err}_H(A) \in \Gamma\left(1, \frac{8}{3}\right), \quad (6)$$

which in turns gives the following probability tail bound

$$\Pr[|\text{err}_H(A)| \geq \epsilon] \leq e^{-\frac{\epsilon^2}{2(1-\frac{8}{3}\epsilon)}}. \quad (7)$$

The proof goes along the following lines: a) we express the error in form of Rademacher chaos b) we use hypercontractive inequalities to bound its moment; this establishes the part (5) c) we plug the moment bound into Taylor’s expansion of the moment generating function and then estimate the series, arriving at the sub-gamma condition (6). By the tail properties of sub-gamma distributions, we conclude the tail bound (7).

2) *Multiple-Sample Estimator:* Next, we move to the multiple-sample estimator in (2). We establish the following

Theorem 2 (Hutchinson’s Estimator n -Sample Bound). *The relative error of n -sample Hutchinson’s Estimator has the following sub-gamma behaviour*

$$\text{err}_H(A) \in \Gamma\left(\frac{1}{n}, \frac{8}{3}\right). \quad (8)$$

In particular, this implies the following bound on the error tail probability, valid for any $0 < \epsilon < 3/8$:

$$\Pr[|\text{err}_H^{(n)}(A)| \geq \epsilon] \leq e^{-\frac{n\epsilon^2}{2(1-\frac{8}{3}\epsilon)}}. \quad (9)$$

This result is obtained from **Theorem 1**, by using extra facts on the concentration of sums of sub-gamma random variables: essentially the variance factor decreases from 1 to $\frac{1}{n}$ for sums of length n (just as the variance of n iid terms).

Corollary 1 (Sample Size for Hutchinson’s Estimator). *The relative error is absolutely bounded by ϵ with probability $1 - \delta$, provided that the sample size n is bigger or equal to*

$$n(\epsilon, \delta) = \frac{2\left(1 - \frac{8}{3}\epsilon\right) \log\left(\frac{1}{\delta}\right)}{\epsilon^2}. \quad (10)$$

This holds for any $0 < \delta < 1$ and $0 < \epsilon < 3/8$.

D. Numerical Comparison with Related Work

Below we demonstrate numerical improvements with respect to the prior works [3], [20]. The bounds, with exact constants, are summarized in Table I below.

Author	Sample Size $n = n(\epsilon, \delta)$
this work	$\frac{2(1-\frac{8}{3}\epsilon)\log(\frac{1}{\delta})}{\epsilon^2}$
[20]	$\frac{6\log(\frac{2}{\delta})}{\epsilon^2}$
[3]	$\frac{6\log(\frac{2}{\delta})+6\text{rank}(A)}{\epsilon^2}$

TABLE I
BOUNDS ON THE ACCURACY OF HUTCHINSON'S ESTIMATOR.

Furthermore, in Figure 2, we present a detailed evaluation of how the sample size n depends on the error ϵ , when the confidence parameter δ is fixed. The setup for this experiment is the discussed hessian of a toy classifier on MNIST.

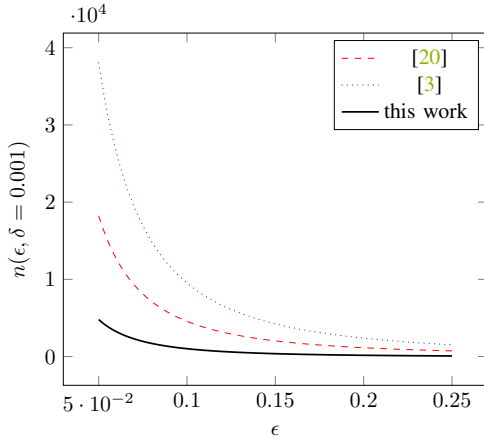


Fig. 2. Sample size for Hutchinson's Estimator. We assume $\delta = 10^{-3}$ which results in confidence of 0.999. The matrix size is 7840, which corresponds to the Hessian of the linear MNIST classification model (input data are black-white images in the resolution of 28×28 , grouped in 10 classes).

III. AUXILIARY RESULTS

1) *Matrix Theory*: We will need the following fact on the decomposition of symmetric matrices [26]. Recall that a matrix B is orthonormal when $B^T = B^{-1}$ (\cdot^T denotes transposition); in particular, the columns and rows of B are of unit length.

Lemma 1 (Symmetric Matrix Factorization). *Any symmetric real matrix A can be written as $A = B^T \Lambda B$ where B is an orthonormal real matrix and Λ is diagonal, consisting of (necessarily real) eigenvalues of A .*

2) *Quadratic Chaos*: We will need the following special case of the celebrated Hypercontractivity Inequality [22], [23]:

Lemma 2 ((2, d)-Hypercontractivity). *Let F be a polynomial of degree 2 in Rademacher variables Z_i (e.g. $F = \sum_{i \neq j} a_{i,j} Z_i Z_j$ for some weights $a_{i,j}$). Then*

$$\|F\|_d \leq (d-1)\|F\|_2. \quad (11)$$

3) *Sub-Gamma Random Variables*: The tail probability of sub-gamma distributions can be estimated from the MGF bound by the Cramer-Chernoff method, as follows:

Lemma 3 (Sub-Gamma Tails). *Let $X \in \Gamma(v, c)$, then*

$$\forall \epsilon > 0 : \Pr[|X| \geq \epsilon] \leq e^{-\frac{\epsilon^2}{v+c\epsilon}}. \quad (12)$$

Moreover, for sub-gamma distributions the following composition property holds for sums of independent terms:

Lemma 4 (Sub-Gamma Sums). *Let random variables $X_i \in \Gamma(v_i, c_i)$ be independent, then*

$$\sum_i X_i \in \Gamma(v', c'), \quad v' \triangleq \sum_i v_i, \quad c' \triangleq \max_i c_i. \quad (13)$$

The proofs of these lemmas can be found in [24].

IV. PROOFS

A. Proof of Theorem 1

Define the following random variable

$$Y = z^T A z - \mathbb{E}[z^T A z] = z^T A z - \text{tr}(A). \quad (14)$$

then our task is to bound the concentration of Y .

1) *Reduction to Off-Diagonal Quadratic Chaos*: Since A is symmetric, by Lemma 1 we have $A = B^T \Lambda B$ where B is orthonormal and Λ is diagonal made of the eigenvalues of A . Since, in addition, the matrix A is positive semi-definite, its eigenvalues are non-negative. Thus, we can write

$$z^T A z = (Bz)^T \Lambda (Bz) = \|\Lambda^{1/2} Bz\|_2^2. \quad (15)$$

This can be decomposed as

$$z^T A z = \sum_i Y_i, \quad Y_i \triangleq (\Lambda^{1/2} Bz)_i^2. \quad (16)$$

To further simplify, denote by $\lambda_1, \dots, \lambda_m$ all the eigenvalues of A ; these are precisely the diagonal entries of Λ . We note that $Y_i = \lambda_i (\sum_j B_{i,j} z_j)^2$, and $\mathbb{E}Y_i = \lambda_i \sum_j B_{i,j}^2$ (because $\mathbb{E}z_j^2 = 1$). Therefore we obtain

$$Y_i - \mathbb{E}Y_i = \lambda_i \sum_{j \neq j'} B_{i,j} B_{i,j'} z_j z_{j'}. \quad (17)$$

Thus, every Y_i is a quadratic form in z_j with zero-diagonal.

2) *Bounding Quadratic Chaos*: For every fixed i we apply Lemma 2 to $Y_i - \mathbb{E}Y_i$, which is a polynomial of degree 2 in Rademacher random variables z_i (note that the off-diagonal property guarantees the degree is exactly two).

Since the random variables $z_j z_{j'}$ indexed by tuples (j, j') for $j \neq j'$ are uncorrelated, we have $\|Y_i - \mathbb{E}Y_i\|_2^2 = \lambda_i^2 \sum_{j \neq j'} B_{i,j}^2 B_{i,j'}^2$; but B is orthonormal, so $\sum_{i,j} B_{i,j}^2 = 1$ and $\|Y_i - \mathbb{E}Y_i\|_2^2 \leq \lambda_i^2$. Thus, we obtain

$$\|Y_i - \mathbb{E}Y_i\|_d \leq \lambda_i^2 (d-1). \quad (18)$$

By the triangle inequality

$$\left\| \sum_i (Y_i - \mathbb{E}Y_i) \right\|_d \leq (d-1) \sum_i \lambda_i. \quad (19)$$

Finally by the definition of Y_i and the standard identity $\text{tr}(A) = \sum_i \lambda_i$ from linear algebra [26]

$$\|z^T Az - \text{tr}(A)\|_d \leq (d-1)\text{tr}(A). \quad (20)$$

Dividing the both sides of this inequality by $\text{tr}(A)$, we complete the proof of the first part of **Theorem 1**.

3) *Bounding Moment Generating Function*: Let $E = z^T Az / \text{tr}(A) - 1$ be the relative error of the estimator. Then by the previous result, the Taylor expansion $e^{tx} = \sum_{d \geq 0} (tx)^d / d!$, and the fact that $\mathbb{E}E = 0$ we have

$$\text{MGF}(E) \leq 1 + \sum_{d \geq 2} \frac{t^d (d-1)^d}{d!}. \quad (21)$$

Let $a_d = (d-1)^d / d!$. Then we have

$$\frac{a_{d+1}}{a_d} = \frac{d^{d+1}}{(d-1)^d} \cdot \frac{1}{d+1}. \quad (22)$$

This ratio is decreasing when $d \geq 2$, and thus

$$\frac{a_{d+1}}{a_d} \leq \frac{8}{3}. \quad (23)$$

Therefore, it follows that

$$\text{MGF}(E) \leq 1 + \frac{t^2}{2(1 - \frac{8}{3}t)}, \quad (24)$$

which, by the standard inequality $\log(1+u) \leq u$, implies

$$\log \text{MGF}(E) \leq \frac{t^2}{2(1 - \frac{8}{3}t)}, \quad (25)$$

so that, by definition, $E \in \Gamma(1, 8/3)$. This completes the proof of the second part of **Theorem 1**.

4) *Concluding Concentration Properties*: The third part of **Theorem 1** follows directly by **Lemma 3**.

B. Proof of **Theorem 2**

We first apply the result on sum of independent sub-gamma distributions from **Lemma 4**, which proves the first part of the theorem. The second part follows by the result on sub-gamma tails, stated in **Lemma 3**.

C. Proof of **Corollary 1**

The corollary follows by equating the right-hand side of the tail bound from **Theorem 2**, and equating with δ . Solving with respect to n gives the claimed formula on the sample size.

V. CONCLUSION

This work establishes a new state-of-art bound on the classical trace estimator due to Hutchinson. The core idea is the clever usage of bounds on Rademacher Chaos and the theory of sub-gamma distributions. The results are numerically superior and of immediate interest to any problems involving trace estimation. Besides, the author hopes to contribute to better popularization of the discussed techniques, novel in the problem context, within the computer science community.

REFERENCES

- [1] M. F. Hutchinson, "A stochastic estimator of the trace of the influence matrix for laplacian smoothing splines," *Communications in Statistics-Simulation and Computation*, vol. 18, no. 3, pp. 1059–1076, 1989.
- [2] Z. Bai, G. Fahey, and G. Golub, "Some large-scale matrix computation problems," *Journal of Computational and Applied Mathematics*, vol. 74, no. 1-2, pp. 71–89, 1996.
- [3] H. Avron and S. Toledo, "Randomized algorithms for estimating the trace of an implicit symmetric positive semi-definite matrix," *J. ACM*, vol. 58, no. 2, Apr. 2011. [Online]. Available: <https://doi.org/10.1145/1944345.1944349>
- [4] L. Lin, Y. Saad, and C. Yang, "Approximating spectral densities of large matrices," *SIAM review*, vol. 58, no. 1, pp. 34–65, 2016.
- [5] I. Han, D. Malioutov, H. Avron, and J. Shin, "Approximating spectral sums of large-scale matrices using stochastic chebyshev approximations," *SIAM Journal on Scientific Computing*, vol. 39, no. 4, pp. A1558–A1585, 2017.
- [6] E. Di Napoli, E. Polizzi, and Y. Saad, "Efficient estimation of eigenvalue counts in an interval," *Numerical Linear Algebra with Applications*, vol. 23, no. 4, pp. 674–692, 2016.
- [7] E. Haber, M. Chung, and F. Herrmann, "An effective method for parameter estimation with pde constraints with multiple right-hand sides," *SIAM Journal on Optimization*, vol. 22, no. 3, pp. 739–757, 2012.
- [8] K. van den Doel and U. Ascher, "Adaptive and stochastic algorithms for eit and dc resistivity problems with piecewise constant solutions and many measurements," *SIAM J. Scient. Comput.*, vol. 34, p. 29, 2012.
- [9] J. Young and D. Ridzal, "An application of random projection to parameter estimation in partial differential equations," *SIAM Journal on Scientific Computing*, vol. 34, no. 4, pp. A2344–A2365, 2012.
- [10] F. Roosta-Khorasani, K. Van Den Doel, and U. Ascher, "Stochastic algorithms for inverse problems involving pdes and many measurements," *SIAM Journal on Scientific Computing*, vol. 36, no. 5, pp. S3–S22, 2014.
- [11] T. van Leeuwen, A. Y. Aravkin, and F. J. Herrmann, "Seismic waveform inversion by stochastic optimization," *International Journal of Geophysics*, vol. 2011, 2011.
- [12] G. H. Golub and U. Von Matt, "Generalized cross-validation for large-scale problems," *Journal of Computational and Graphical Statistics*, vol. 6, no. 1, pp. 1–34, 1997.
- [13] H. Avron, "Counting triangles in large graphs using randomized matrix trace estimation," in *Workshop on Large-scale Data Mining: Theory and Applications*, vol. 10, 2010, pp. 10–9.
- [14] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [15] B. Christianson, "Automatic hessians by reverse accumulation," *IMA Journal of Numerical Analysis*, vol. 12, no. 2, pp. 135–150, 1992.
- [16] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard et al., "Tensorflow: A system for large-scale machine learning," in *12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16)*, 2016, pp. 265–283.
- [17] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga et al., "Pytorch: An imperative style, high-performance deep learning library," in *Advances in neural information processing systems*, 2019, pp. 8026–8037.
- [18] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang, "JAX: composable transformations of Python+NumPy programs," 2018. [Online]. Available: <http://github.com/google/jax>
- [19] S. Ubaru, J. Chen, and Y. Saad, "Fast estimation of $\text{tr}(f(a))$ via stochastic lanczos quadrature," *SIAM Journal on Matrix Analysis and Applications*, vol. 38, no. 4, pp. 1075–1099, 2017.
- [20] F. Roosta-Khorasani and U. Ascher, "Improved bounds on sample size for implicit matrix trace estimators," *Foundations of Computational Mathematics*, vol. 15, no. 5, pp. 1187–1212, 2015.
- [21] Z. Yao, A. Gholami, K. Keutzer, and M. Mahoney, "Pyhessian: Neural networks through the lens of the hessian," *arXiv preprint arXiv:1912.07145*, 2019.
- [22] R. O'Donnell, *Analysis of boolean functions*. Cambridge University Press, 2014.

- [23] A. Bonami, “Ensembles $\lambda(p)$ dans le dual de d^∞ ,” *Annales de l’institut Fourier*, vol. 18, no. 2, pp. 193–204, 1968. [Online]. Available: <http://eudml.org/doc/73956>
- [24] S. Boucheron, G. Lugosi, and P. Massart, *Concentration Inequalities: A Nonasymptotic Theory of Independence*. OUP Oxford, 2013. [Online]. Available: <https://books.google.at/books?id=koNqWRlulhPOC>
- [25] R. Vershynin, *High-dimensional probability: An introduction with applications in data science*. Cambridge university press, 2018, vol. 47.
- [26] D. Mello, *Invitation to Linear Algebra*, ser. Invitation to Linear Algebra. CRC Press, 2017, no. v. 978, nos. 1-7952. [Online]. Available: <https://books.google.at/books?id=wEI-MQAACAAJ>